



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

BSC tools to study the computational efficiency of EC-Earth components

Miguel Castrillo, Oriol Tintó, Kim Serradell





Objectives

- Earth System Model
- Reliable in-house predictions of global climate change
- Part of a Europe-wide consortium
- Being used in large European projects
 - EMBRACE
 - EUPORIAS
 - IS-ENES
 - SPECS
- 3.1 version → IFS + NEMO-LIM + OASIS



- Energy efficiency

$$EE = \frac{\text{Performance}}{\text{Power consumed}}$$

← Increase performance

← Reduce power

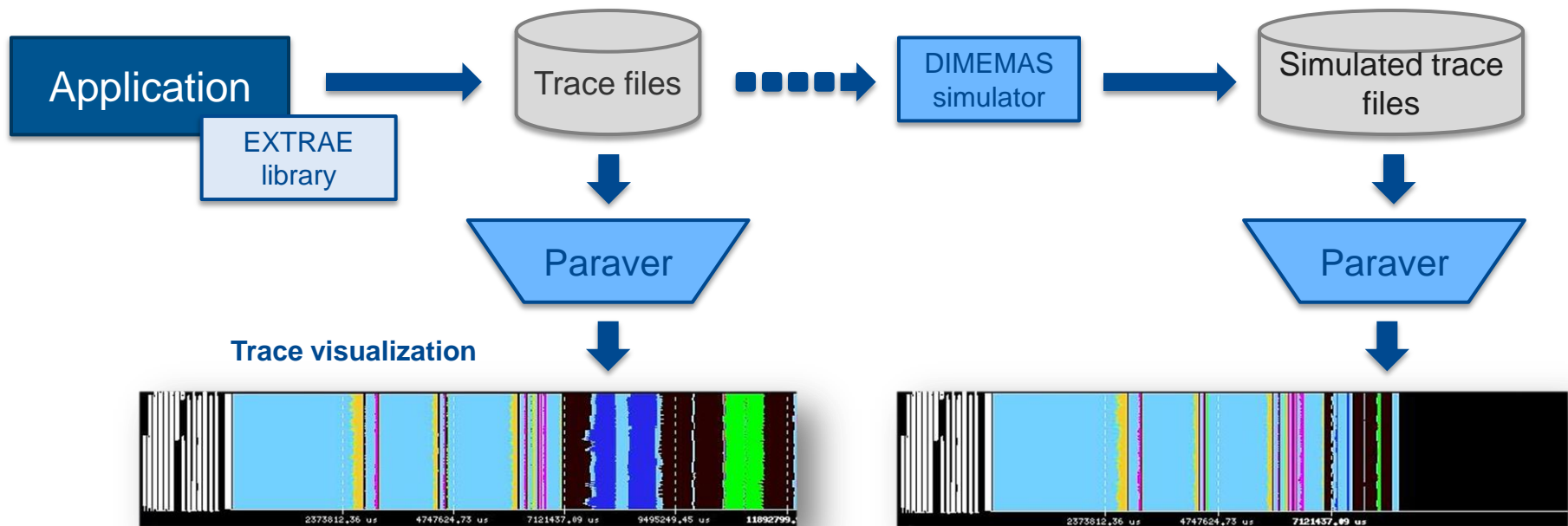
- Performance loss caused by:

- Bad programming
- Load imbalance
- Synchronization
- Resource contention
- ...



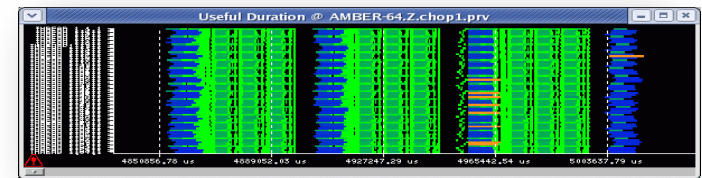
Methodology

- Since 1991
- Based on traces
- Open Source: <http://www.bsc.es/paraver>
- **Extrae**: Package that generates Paraver trace-files for a post-mortem analysis
- **Paraver**: Trace visualization and analysis browser
 - Includes trace manipulation: Filter, cut traces
- **Dimemas**: Message passing simulator



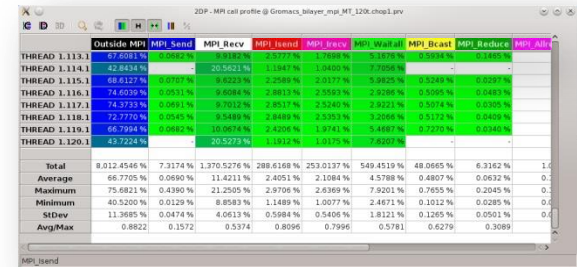
- BSC instrumentation package
- When/Where
 - Parallel programming model runtime
 - Selected user functions
 - Periodic samples
 - User events
- Additional information
 - Counters

- Every behavioral aspect/metric of a thread can be described as a function of time
- Those functions of time can be rendered into a 2D image



- Statistics can be computed for each possible value or range of values of that function of time

- Tables: Profiles and histograms



	Outside MPI	MPI Send	MPI Recv	MPI Isend	MPI Rcvv	MPI Waitall	MPI Bcast	MPI Reduce	MPI All
THREAD 1.113.1	67.8501%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.114.1	62.8835%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.115.1	58.6127%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.116.1	74.6036%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.117.1	74.3733%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.118.1	72.7702%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.119.1	66.7964%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
THREAD 1.120.1	43.7224%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%	0.0000%
Total	8.0124546%	7.3174%	1.3705276%	288.6168%	253.0137%	569.4519%	48.0665%	6.3362%	1.0
Average	66.7705%	0.0690%	11.4211%	2.4051%	2.1084%	4.5788%	0.4807%	0.0632%	0.0
Maximum	75.6821%	0.4390%	21.2505%	2.8708%	2.6369%	7.9201%	0.7655%	0.2045%	0.0
Minimum	40.5200%	0.0129%	8.8583%	1.1489%	1.0077%	2.4671%	0.1012%	0.0285%	0.0
StdDev	11.3685%	0.0414%	4.6613%	0.5684%	0.5405%	1.8121%	0.1265%	0.0501%	0.0
Avg/Max	0.8821	0.1572	0.5374	0.8096	0.7996	0.5781	0.6279	0.3089	

- Types of functions

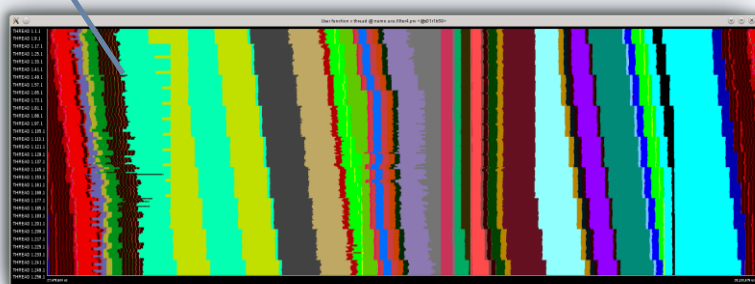
- Categorical: State, user function...
 - Logical: In specific user function, in MPI call...
 - Numerical: IPC, L2 miss ratio, duration of computation burst...

PARAVER trace analysis

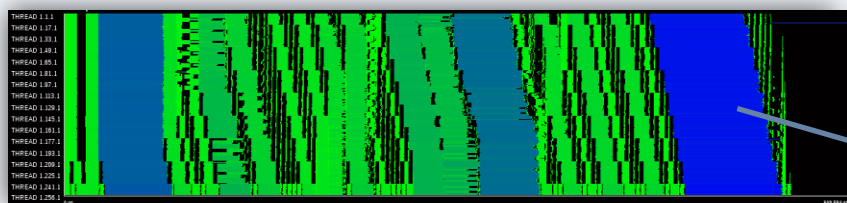
Serial efficiency

Each color represents a function

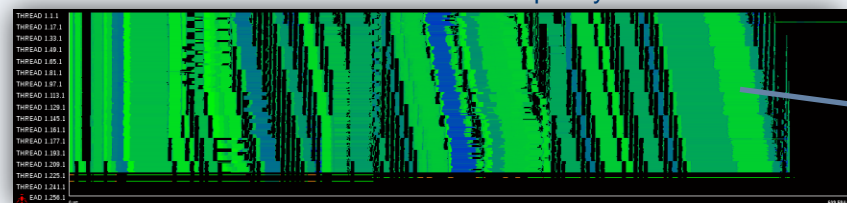
Functions instrumentation



Useful Duration



Useful IPC - Instructions per cycle



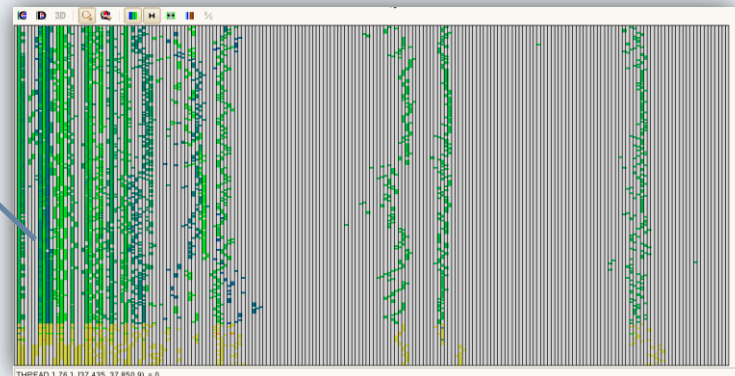
Horizontal axis -> Time component

Darker color =
Higher
values

Darker color =
Longer
duration

Lighter color =
Less
IPC

Correlation between two functions



Parallel efficiency

	Outside MPI	MPI_Allrecv	MPI_Isend	MPI_Wait	MPI_Allreduce	MPI_Allgather
THREAD 1.236.1	70.43 %	28.58 %	0.12 %	0.18 %	0.09 %	
THREAD 1.237.1	71.34 %	26.58 %	0.18 %	1.23 %	0.75 %	
THREAD 1.238.1	72.11 %	26.53 %	0.11 %	0.54 %	0.72 %	
THREAD 1.239.1	71.48 %	26.94 %	0.10 %	0.79 %	0.68 %	
THREAD 1.240.1	70.58 %	28.08 %	0.11 %	0.81 %	0.68 %	
THREAD 1.241.1	69.54 %	29.80 %	0.18 %	0.18 %	0.18 %	0.82 %
THREAD 1.242.1	68.36 %	3.80 %	0.18 %	0.68 %	0.05 %	0.90 %
THREAD 1.243.1	69.08 %	3.71 %	0.13 %	0.40 %	0.05 %	0.90 %
THREAD 1.244.1	68.78 %	0.98 %	0.15 %	0.17 %	0.05 %	0.90 %
THREAD 1.245.1	68.91 %	1.01 %	0.14 %	0.12 %	0.05 %	0.97 %
THREAD 1.246.1	68.74 %	2.82 %	0.14 %	0.14 %	0.05 %	0.91 %
THREAD 1.247.1	68.90 %	1.70 %	0.18 %	0.21 %	0.05 %	0.91 %
THREAD 1.248.1	68.32 %	1.38 %	0.17 %	0.18 %	0.05 %	0.90 %
THREAD 1.249.1	68.14 %	1.43 %	0.13 %	0.20 %	0.07 %	10.93 %
THREAD 1.250.1	68.24 %	1.28 %	0.12 %	0.59 %	0.07 %	9.72 %
THREAD 1.251.1	67.00 %	2.31 %	0.18 %	0.10 %	0.07 %	10.34 %
THREAD 1.252.1	67.58 %	1.80 %	0.13 %	0.15 %	0.08 %	10.13 %
THREAD 1.253.1	67.35 %	1.88 %	0.18 %	0.09 %	0.08 %	10.30 %
THREAD 1.254.1	67.71 %	1.47 %	0.14 %	0.40 %	0.08 %	10.21 %
THREAD 1.255.1	67.73 %	1.80 %	0.18 %	1.38 %	0.08 %	9.80 %
THREAD 1.256.1	67.66 %	0.93 %	0.17 %	0.21 %	0.08 %	9.42 %
Total	18,718.18 %	4,259.58 %	61.57 %	2,133.74 %	276.02 %	150.84 %
Average	73.12 %	16.64 %	0.24 %	8.33 %	1.08 %	9.43 %
Maximum	89.58 %	28.60 %	0.31 %	23.62 %	1.30 %	10.34 %
Minimum	69.38 %	0.82 %	0.10 %	0.06 %	0.04 %	0.26 %
StdDev	4.00 %	8.68 %	0.04 %	8.11 %	0.32 %	0.89 %
AvgMax	0.82	0.58	0.77	0.35	0.83	0.91

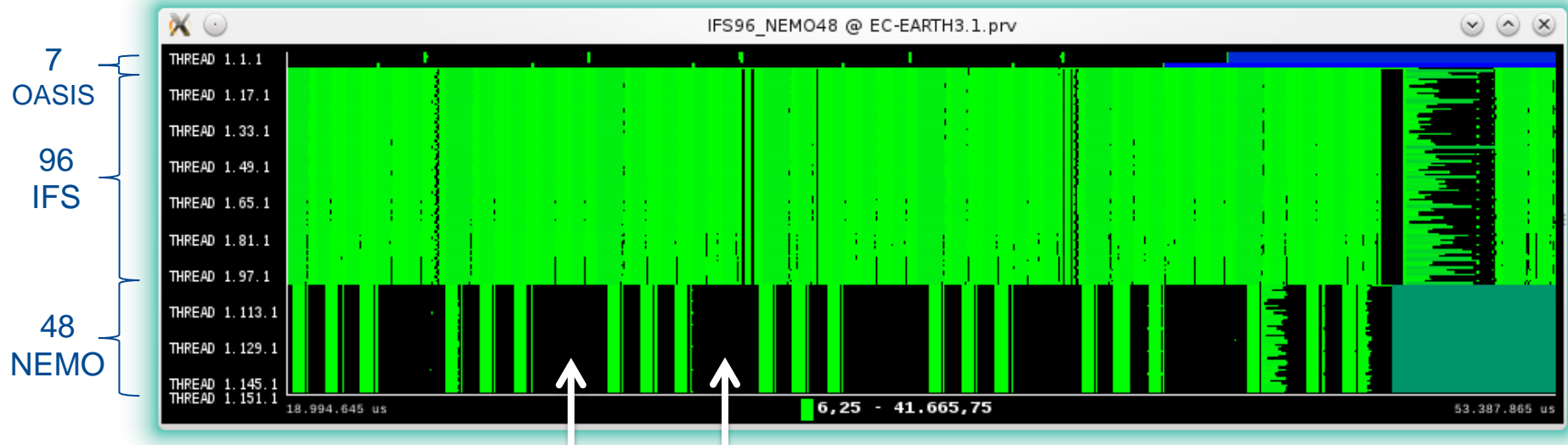
MPI Stats reflect the percentage of time invested in computation for each thread.

Total stats give the communication efficiency and the load balance

An EC-Earth Paraver trace



- Motivation: Finding a good configuration to optimize the resources usage.
- IFS T255L91-ORCA1L46
- Configuration used in production
 - Using 7 cores for OASIS, 96 for IFS and 48 for NEMO
- 1 day simulation traces
- Traces generated in burst mode (only computational regions > 100us)
- Paraver view → Useful duration (displays duration of computational bursts)



Black regions → Not computation (MPI, I/O...) → NEMO waiting

Time axis

Dealing with resource contention



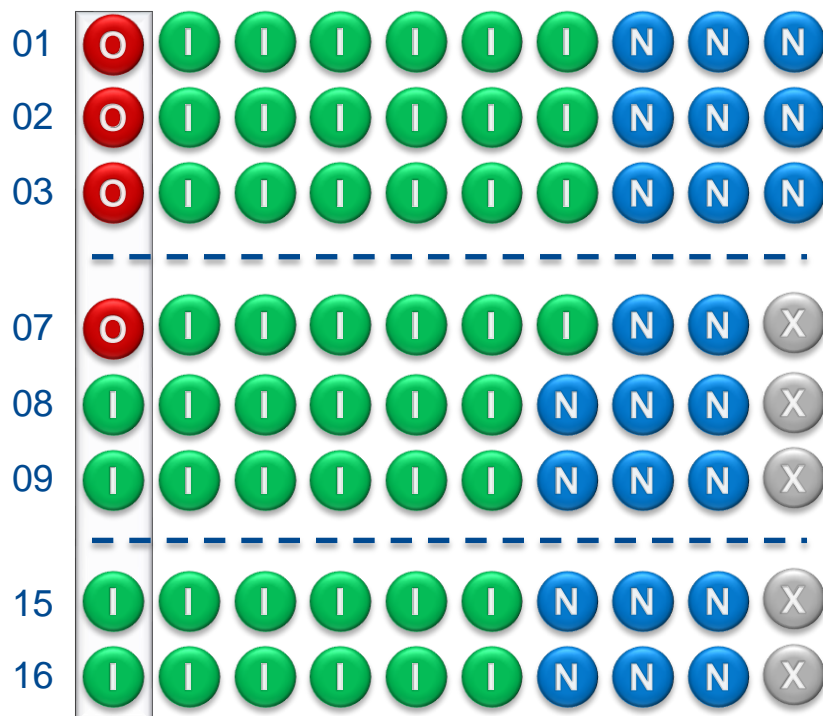
- Resource contention: Conflict over access to shared resources, i.e. memory, storage, buses...
- Distributing the processes among different nodes to mitigate the problem.

 7 Oasis +  96 IFS +  48 NEMO in 10 nodes (16 cores per node)

Default Mode

Nodes

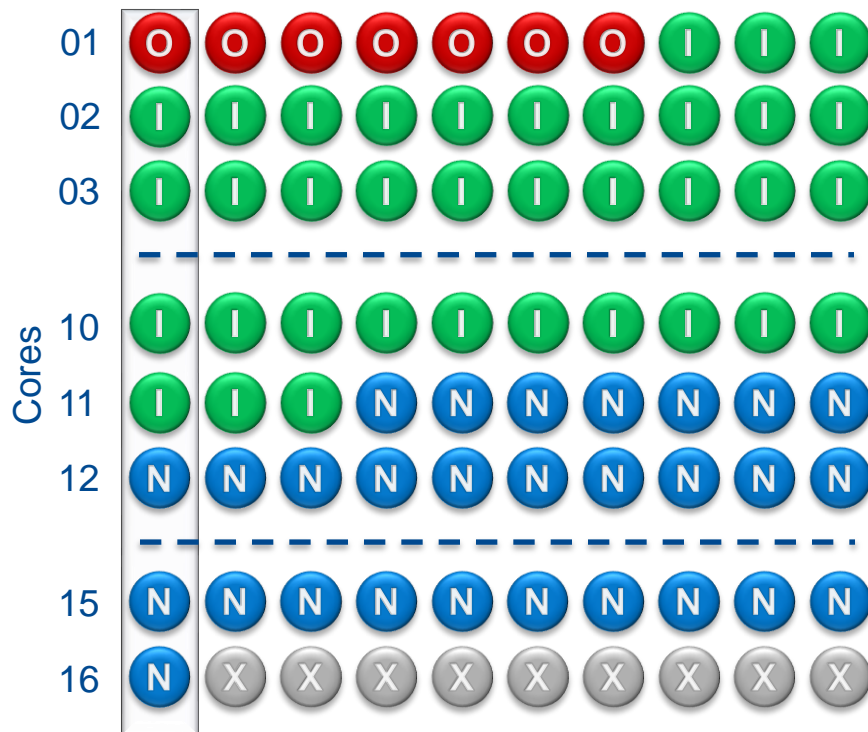
01 02 03 04 05 06 07 08 09 10



Processor Affinity

Nodes

01 02 03 04 05 06 07 08 09 10

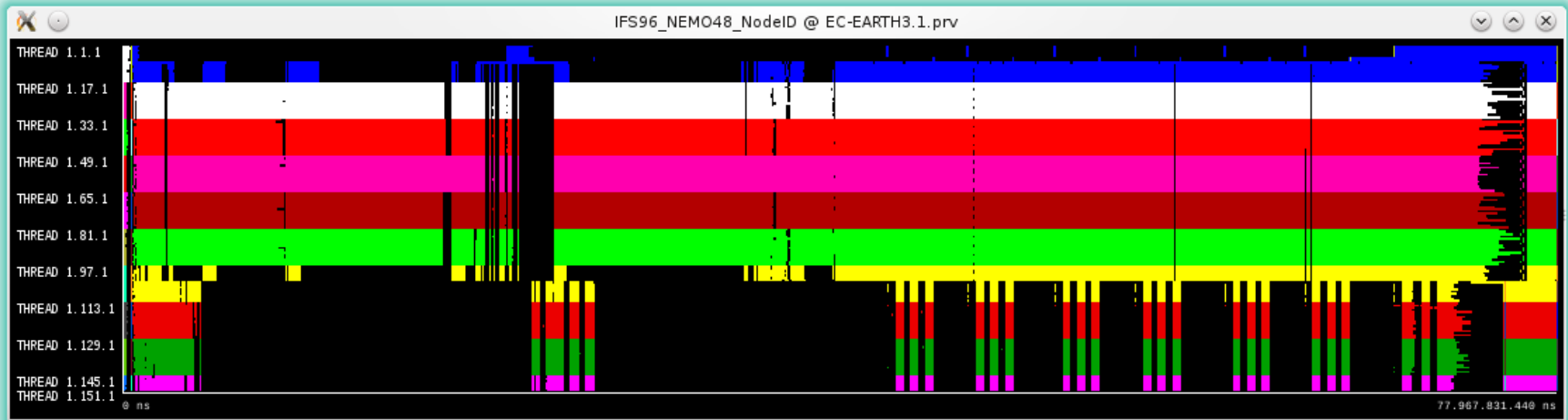


Dealing with resource contention

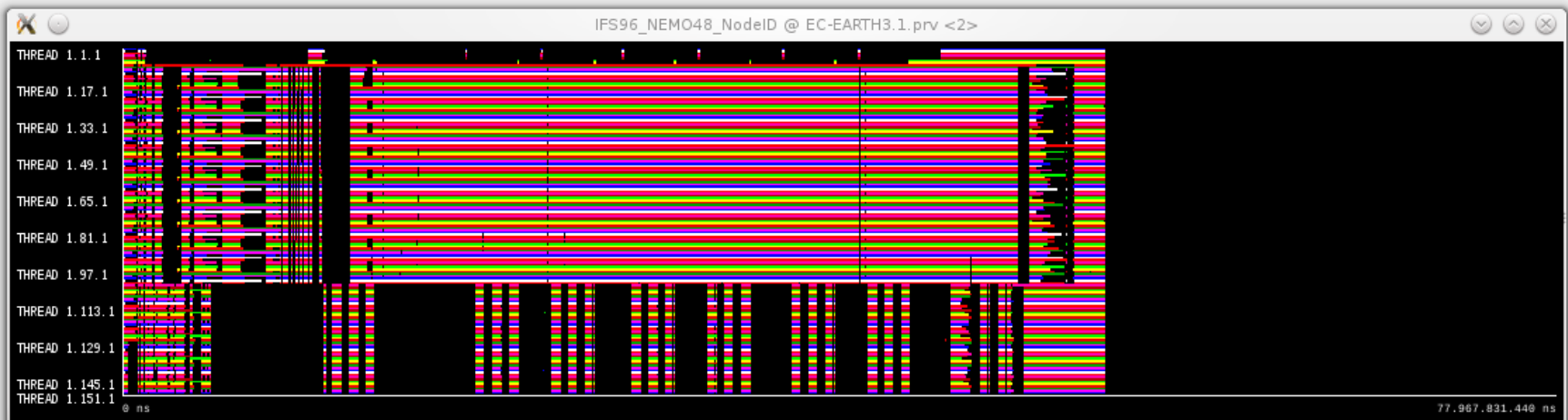


- Color identifies the different nodes

Default Mode

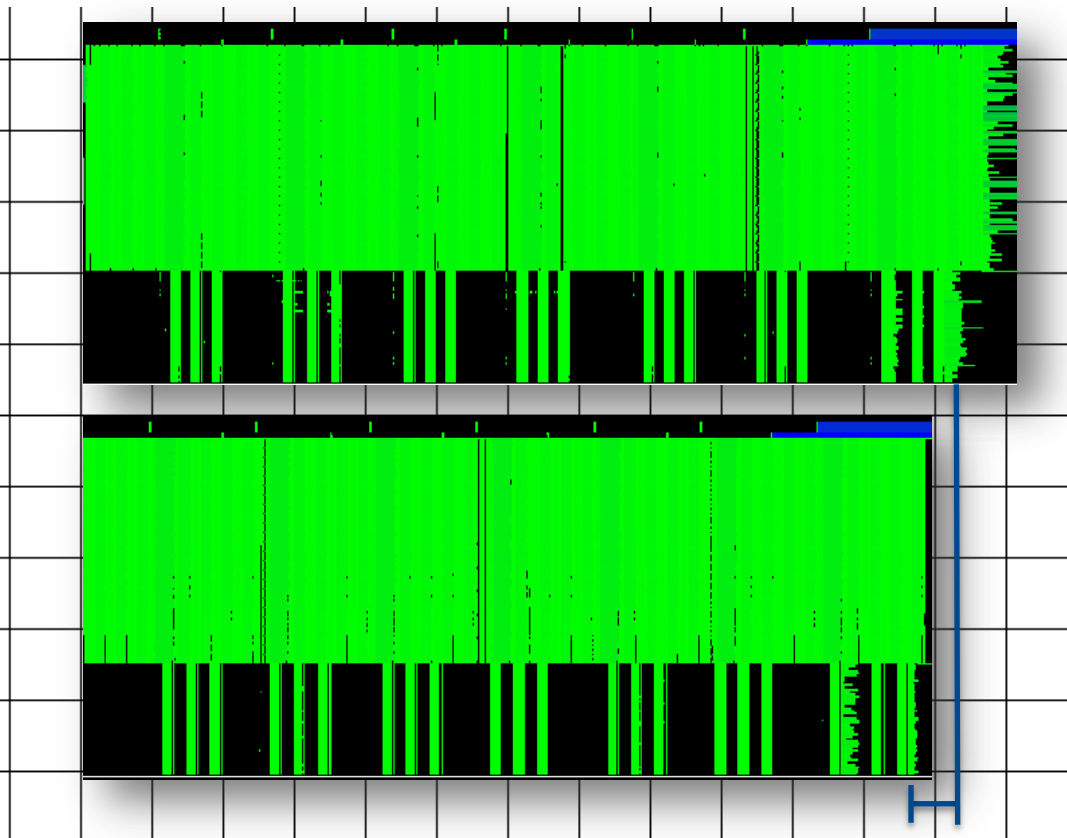


Processor Affinity



- There is an improvement, but NEMO is still waiting a lot for IFS
- Next step: Increase IFS and reduce NEMO resources to find an equilibrium

32,76 s



2,1 s

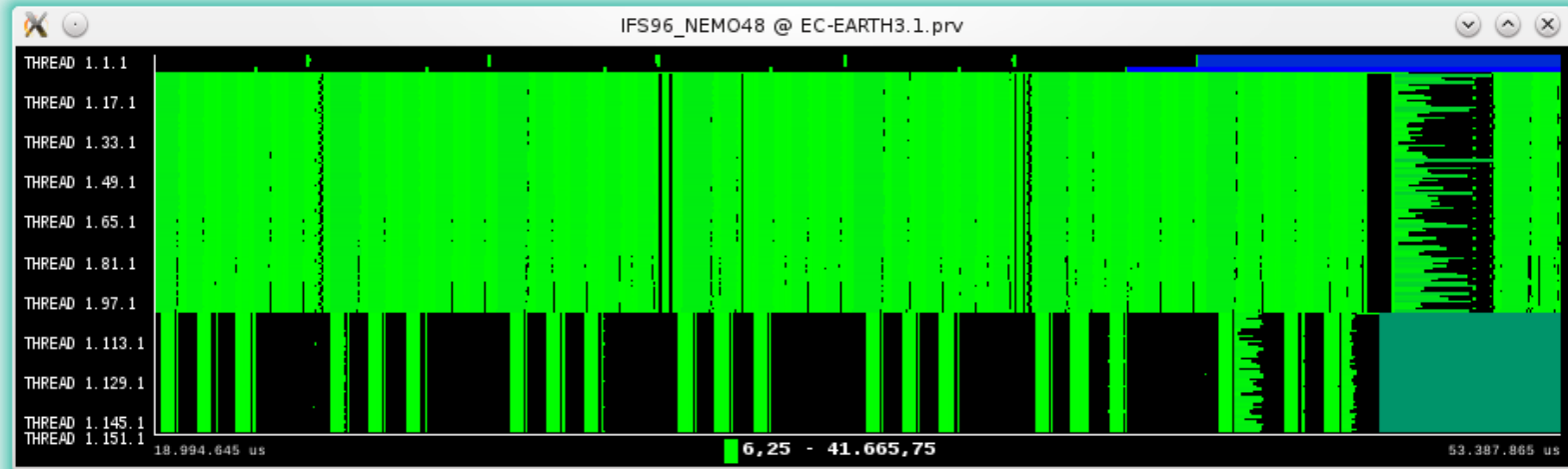
1 day simulation → 2,1s save
10 years simulation → 2,1h save

Only with a change in the
order of processes !!

Finding an optimum configuration



- Finding a better configuration to increase the performance.

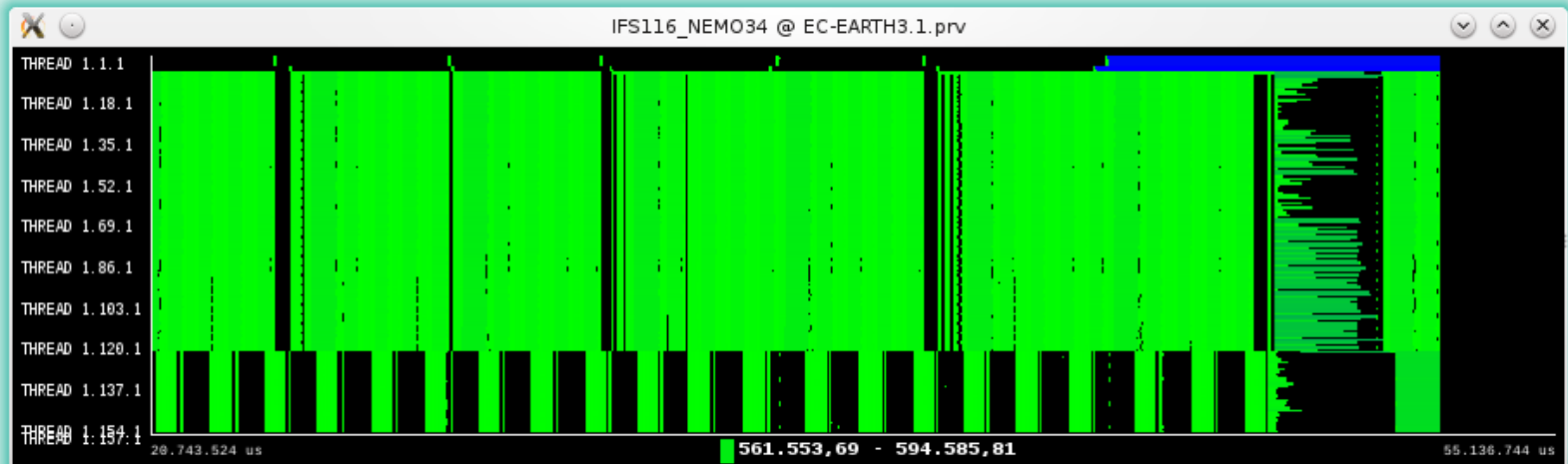
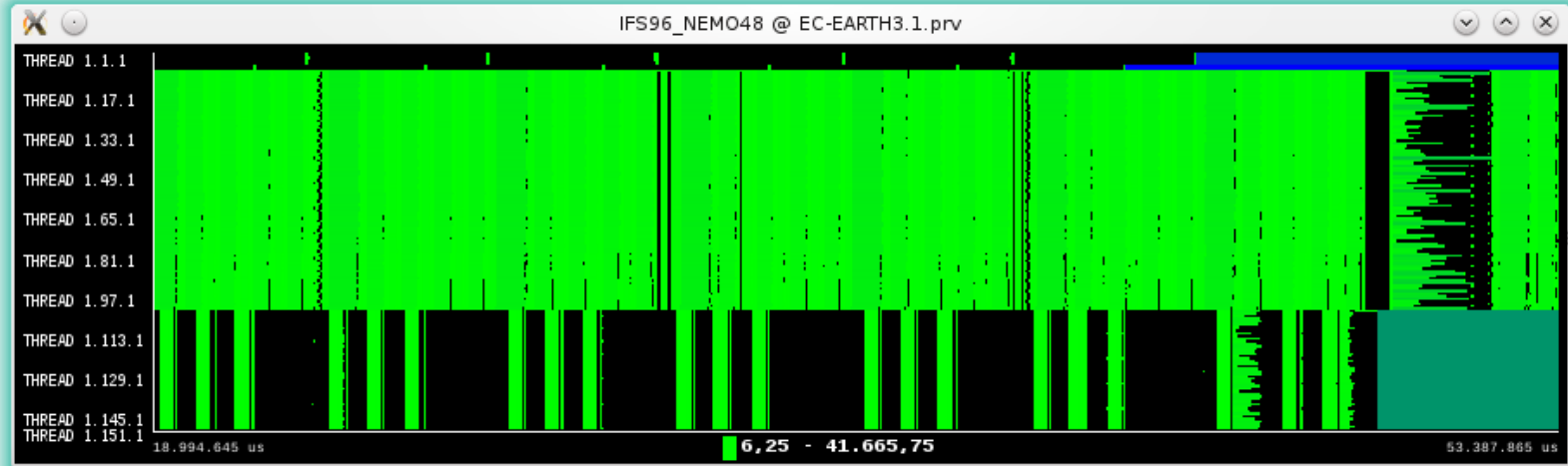


2300 ms waiting

Finding an optimum configuration



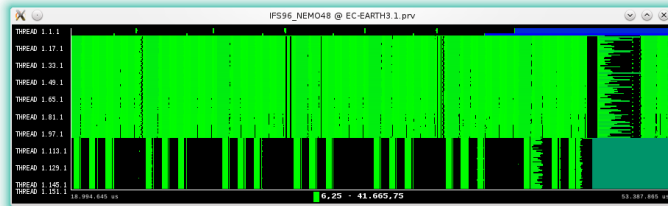
- Finding a better configuration to increase the performance.
- Optimum configuration for this resolution: 116 IFS – 34 NEMO



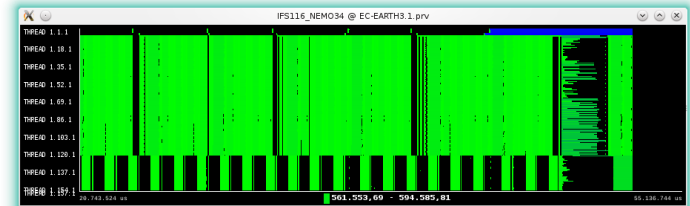
Finding an optimum configuration

- The time gain comes from a better resource usage
- The workload is better balanced → Can be confirmed in Paraver

7 OASIS - 96 IFS – 48 NEMO



7 OASIS – 116 IFS – 34 NEMO



2DP - efficiency @ EC-EART...

	[0..1]
Total	99,90
Average	0,66
Maximum	0,82
Minimum	0,36
StDev	0,17
Avg/Max	0,80

Comms.
Efficiency

Load
Balance

2DP - efficiency @ EC-EART...

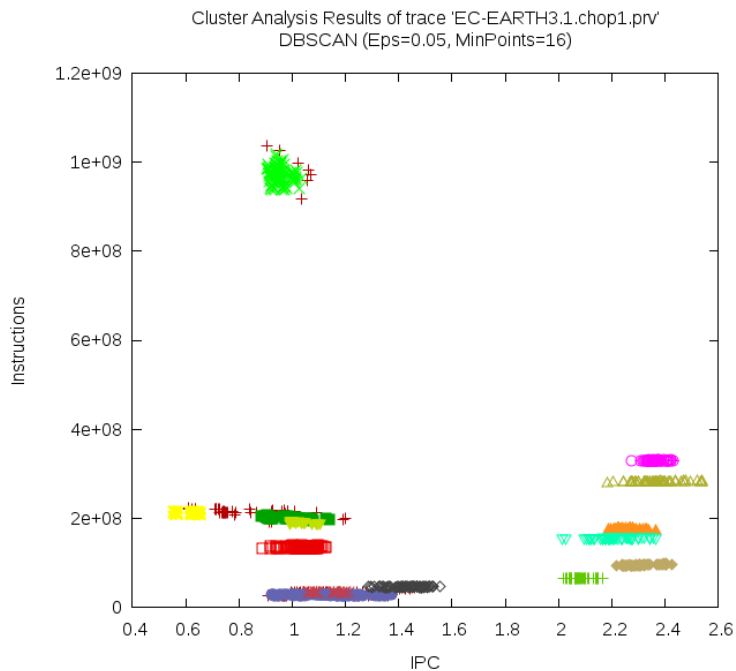
	[0..1]
Total	108,96
Average	0,69
Maximum	0,80
Minimum	0,08
StDev	0,13
Avg/Max	0,87

 **+0.07
LB efficiency**

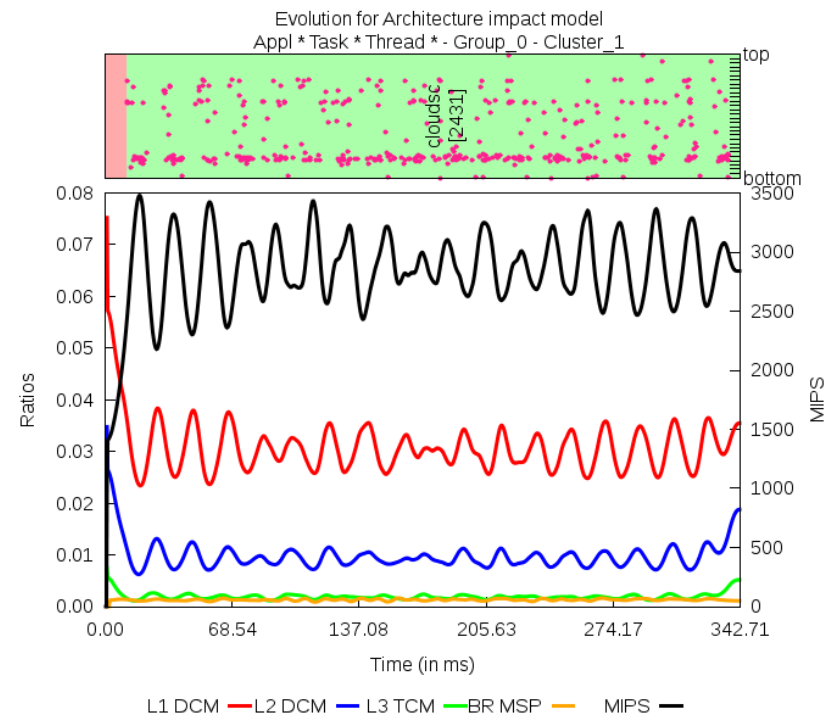


Ongoing work

- **Clustering** → Data mining technique for data classification
 - Detects different trends in the computation regions (CPU bursts) of an application
 - Similarity intended in terms of duration or hardware counter reduced metrics
- **Folding** → Instantaneous metrics with minimum overhead
 - Instrumentation (delimits regions) + sampling (progression within the regions)
 - Captures performance counters and call-stack references

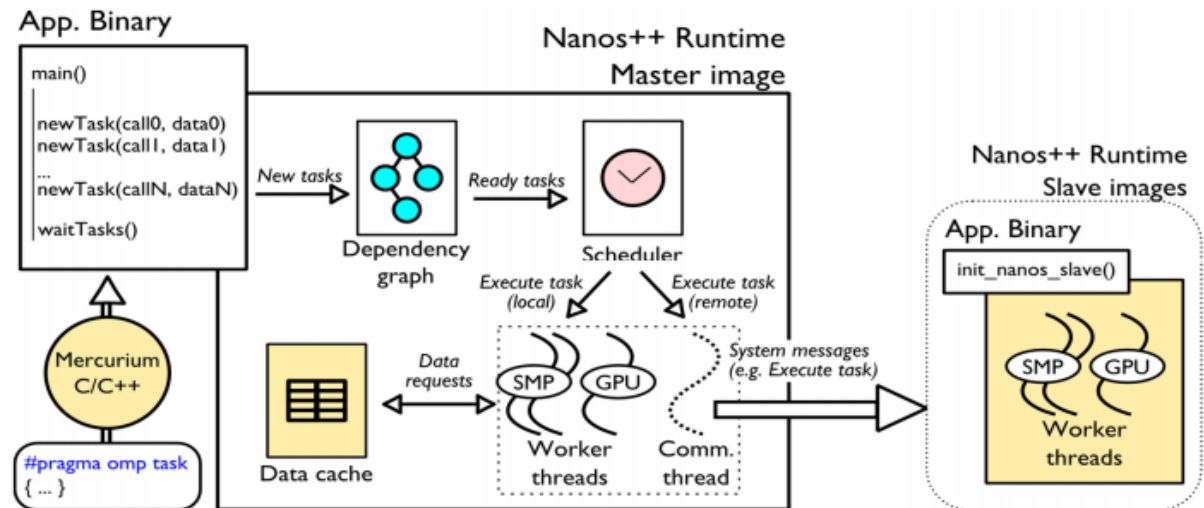


EC-Earth clustering scatter plot

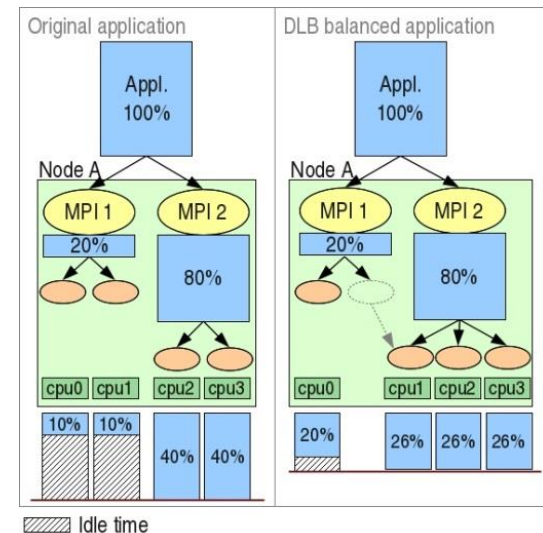
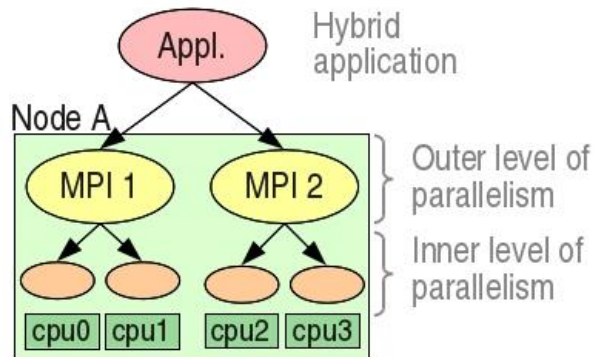


Performance view of a folded cluster

- Parallel Programming Model
- Build on existing standard: OpenMP
- Directive based to keep a serial version
- Targeting: SMP, clusters, and accelerator devices
- Developed at the Barcelona Supercomputing Center (BSC)
 - Mercurium source-to-source compiler
 - Nanos++ runtime system
- <https://pm.bsc.es/ompss>

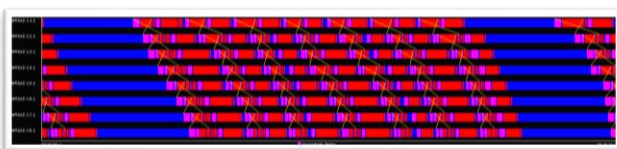


- Dynamic library
- Speeds up hybrid applications with nested parallelism
 - Improving the load imbalance of the outer level parallelism
 - Different load balancing algorithms
- Load balance within node
- Automatically achieved by the runtime
- LeWI: Lend core When Idle

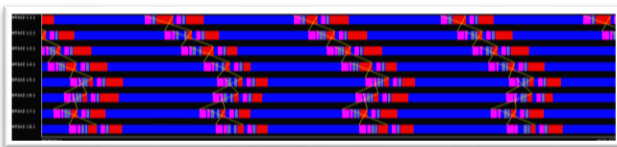


Communication Bottlenecks identified using ORCA1 domain

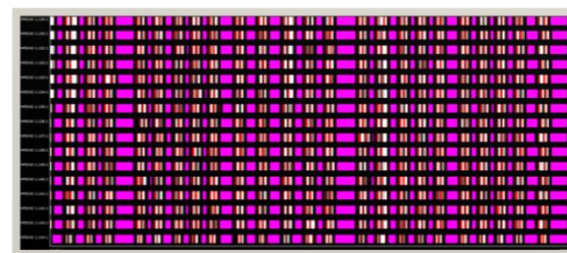
Computation / Communication ratio



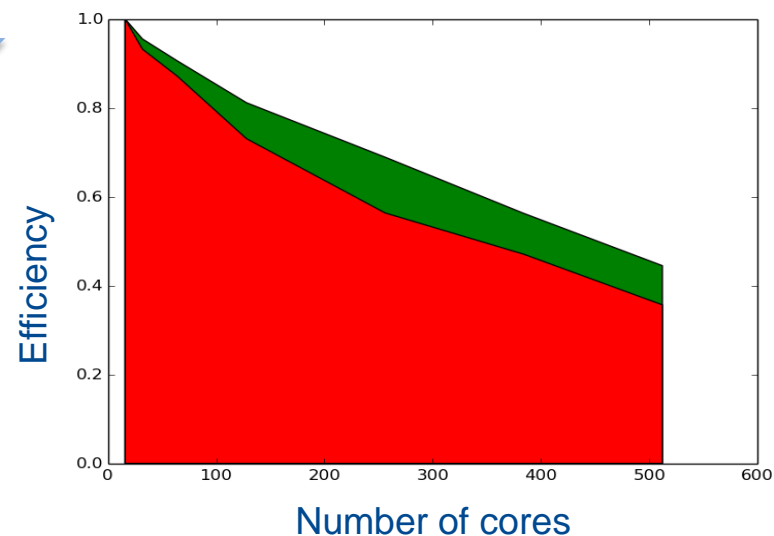
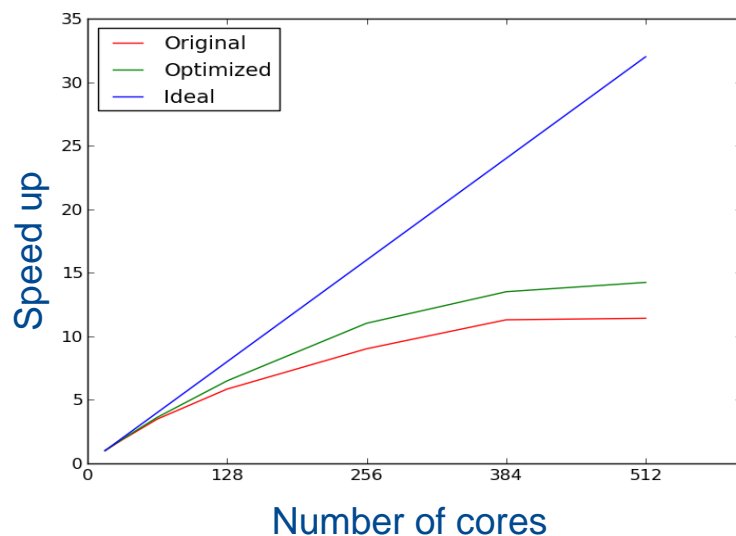
Message Packing



Abuse of Global Communications



Reduction of Global Communications





Conclusions

- Little changes in the configuration can significantly improve the performance.
- Trace analysis can guide the users in performing this task.
- A precise analysis and prediction can generate ideas that direct the restructuring of the application in the most productive way.



**Barcelona
Supercomputing
Center**

Centro Nacional de Supercomputación



EXCELENCIA
SEVERO
OCHOA

Thank you!

For further information please contact
miguel.castrillo@bsc.es